

CHAPTER #9

①

PROBLEM OF THIS CHAPTER: FOR A CONVEX $f(x)$,
minimize: $f(x)$ I.E. NO CONSTRAINTS. ALSO,
 $f(x)$ IS TWICE CONTINUOUSLY DIFFERENTIABLE (\Rightarrow DOM f IS OPEN)

NECESSARY AND SUFFICIENT CONDITION FOR A POINT TO BE
OPTIMAL: $\nabla f(x^*) = 0$, WHICH IS SYSTEM OF n
EQUATIONS WITH n UNKNOWN. HOWEVER, WE CAN HARDLY
SOLVE IT IN PRACTICE. SO WE FIND SEQUENCE
 $x^{(k)}$ WITH $f(x^{(k)}) \rightarrow p^*$ AS $k \rightarrow \infty$
MINIMIZING SEQUENCE AND STOP WHEN $f(x^{(k)}) - p^* \leq \epsilon$

REQUIREMENT: WE NEED A STARTING POINT $x^{(0)} \in \text{dom} f$

SUCH THAT THE SUBLEVEL SET

$$S = \{x \in \text{dom} f \mid f(x) \leq f(x^{(0)})\} \text{ IS CLOSED}$$

(OTHERWISE, THE ABOVE SEQUENCE MAY CONVERGE TO A
POINT NOT IN S). THIS IS GUARANTEED IF f IS CLOSED

DEFINITION: f IS CLOSED IF ALL ITS SUBLEVEL SETS ARE
CLOSED (SUCH AS \mathbb{R}^m)

PROPERTY: 1) CONTINUOUS FUNCTIONS WITH CLOSED DOMAINS ARE CLOSED
2) CONTINUOUS FUNCTIONS WITH OPEN DOMAINS FOR WHICH $f(x)$
TENDS TO INFINITY AS $x \rightarrow \text{bd dom} f$ ARE CLOSED.

EXAMPLE #1 minimize $\frac{1}{2} x^T P x + q^T x + r$,

$$P \in S_+^m, q \in \mathbb{R}^m, r \in \mathbb{R}$$

$$\nabla \left(\frac{1}{2} x^T P x + q^T x + r \right) = 0 \Leftrightarrow \boxed{P x + q = 0} \text{ WHICH}$$

ARE LINEAR

CASES: 1) ONE SOLUTION (ALWAYS THE CASE WHEN $p > 0$) (2)
 \Rightarrow UNIQUE MINIMUM

2) INFINITE SOLUTIONS (INFINITE MINIMA)

3) NO SOLUTIONS \Rightarrow UNBOUNDED BELOW.

NOTE THAT SUBLEVEL SET CONDITION IS SATISFIED

EXAMPLE #2

minimize $f(x) = \log \left(\sum_{i=1}^m \exp[a_i^T x + b_i] \right)$

OPTIMALITY CONDITION

$$\frac{\partial f(x)}{\partial a_j} = 0 \Leftrightarrow$$

$$\frac{1}{\log \left(\sum_{j=1}^m \exp[a_j x + b_j] \right)} \sum_{i=1}^m \exp[a_i^T x + b_i] a_{ij} = 0$$

THEREFORE,

$$\nabla f(x) = 0 \Leftrightarrow \frac{1}{\sum_{j=1}^m \exp[a_j^T x + b_j]} \cdot \sum_{i=1}^m \exp[a_i^T x + b_i] a_i = 0$$

WHICH CANNOT BE SOLVED ANALYTICALLY, SO HENCE WE NEED ITERATIVE METHOD. AGAIN, SUBLEVEL SET CONDITION IS SATISFIED

EXAMPLE #3

minimize $f(x) = - \sum_{i=1}^m \log(b_i - a_i^T x)$,

I.E. THE LOGARITHMIC BARRIER FOR THE SET OF THE INEQUALITIES. AGAIN, THE CONDITION IS SATISFIED. REGARDING THE OPTIMALITY CONDITION,

$$\frac{\partial f}{\partial a_j} = 0 \Leftrightarrow - \sum_{i=1}^m \frac{1}{b_i - a_i^T x} \cdot (-a_{ij}) = \sum_{i=1}^m \frac{a_{ij}}{b_i - a_i^T x} = 0,$$

WHICH AGAIN CANNOT BE SOLVED ANALYTICALLY

IN VECTOR FORM,

$$\nabla f(x) = 0 \Leftrightarrow \sum_{i=1}^m \frac{a_i}{b_i - a_i^T x} = 0$$

AGAIN, SUBLEVEL CONDITION IS SATISFIED

9.1.2 STRONG CONVEXITY

LET f BE TWICE CONTINUOUSLY DIFFERENTIABLE.

DEFINITION: f IS STRONGLY CONVEX ON S IF $f'' \geq mI$.

(9.7) $\nabla^2 f(x) \geq mI \iff \nabla^2 f(x) - mI \succeq 0 \quad \forall x \in S.$

SO LET US ASSUME THAT f IS STRONGLY CONVEX.

PROPERTY 1: $\forall x, y \in S,$
 $f(y) \geq f(x) + \nabla f(x)^T (y-x) + \frac{m}{2} \|y-x\|_2^2$
($m=0 \implies$ KNOWN CASE) (9.8)

PROOF:

KNOWN PROPERTY: FOR $x, y \in S,$

$$f(y) = f(x) + \nabla f(x)^T (y-x) + \frac{1}{2} (y-x)^T \nabla^2 f(z) (y-x)$$

WHERE z ON THE LINE SEGMENT $[x, y]$ BY (1),
LAST TERM IS AT LEAST $\frac{m}{2} \|x-y\|_2^2$, SO RESULT FOLLOWS.

PROPERTY 2: LET ANY $x \in S$. THEN

$$f(x) - p^* \leq \frac{1}{2m} \|\nabla f(x)\|_2^2, \quad \text{SO WE BOUND SUBOPTIMALITY}$$

PROOF:

WE KNOW $f(y) \geq f(x) + \nabla f(x)^T (y-x) + \frac{m}{2} \|y-x\|_2^2$

IT IS MINIMIZE THE RIGHT HAND SIDE W.R.T y :

$$\nabla \left[f(x) + \nabla f(x)^T (y-x) + \frac{m}{2} (y-x)^T I (y-x) \right] = 0 \iff$$

$$\nabla f(x) + \frac{m}{2} \cdot 2 (y-x) = 0 \iff$$

$$\tilde{y} = x - \frac{1}{m} \nabla f(x), \quad \text{THEREFORE}$$

$$\begin{aligned} f(y) &\geq f(x) + \nabla f(x)^T [\tilde{y} - x] + \frac{m}{2} \|\tilde{y} - x\|_2^2 \\ &= f(x) + \nabla f(x)^T \left(-\frac{1}{m} \nabla f(x) \right) + \frac{m}{2} \left(\frac{1}{m^2} \right) \|\nabla f(x)\|_2^2 \\ &= f(x) - \frac{1}{2m} \|\nabla f(x)\|_2^2 \end{aligned}$$

THIS BOUND HOLD FOR ALL y , SO IT WILL ALSO HOLD FOR $\inf_y f(y)$, SO $p^* \geq f(x) - \frac{L}{2m} \|\nabla f(x)\|_2^2$ AND THE RESULT FOLLOWS

PROPERTY 3: $\|\nabla f(x)\|_2 \leq (2m\epsilon)^{\frac{1}{2}} \Rightarrow f(x) - p^* \leq \epsilon$
 (PROOF COMES IMMEDIATELY FROM PREVIOUS.)

PROPERTY 4: $\|x - x^*\|_2 \leq \frac{2}{m} \|\nabla f(x)\|_2$ (THEREFORE x^* IS UNIQUE)

PROOF: WE APPLY (9.8) OF PREVIOUS PAGE FOR $y = x^*$
 $\Rightarrow p^* = f(x^*) \geq f(x) + \nabla f(x)^T (x^* - x) + \frac{m}{2} \|x^* - x\|_2^2$
 $\geq f(x) - \|\nabla f(x)\| \|x^* - x\| + \frac{m}{2} \|x^* - x\|_2^2$
 (BY CAUCHY-SCHWARZ), AND BECAUSE $p^* - f(x) \leq 0$, IT FOLLOWS THAT $-\|\nabla f(x)\| \|x^* - x\| + \frac{m}{2} \|x^* - x\|_2^2 \leq 0$, AND THE RESULT FOLLOWS.

PROPERTY 5: THE HESSIAN IS ALSO UPPER BOUNDED! IN PARTICULAR ANY SUBLEVEL SET IS BOUNDED, AND BECAUSE THE MAXIMUM EIGENVALUE IS A CONTINUOUS FUNCTION OF x , IT IS ALSO UPPER BOUNDED

$$\exists M : \nabla^2 f(x) \leq M I$$

PROPERTY 6:

$$f(y) \leq f(x) + \nabla f(x)^T (y - x) + \frac{M}{2} \|y - x\|_2^2 \quad (9.13)$$

Proof: R.H.S IS $f(x) + \nabla f(x)^T (y-x) + \frac{M}{2} (y-x)^T I (y-x)$ (5)

SET $z = y-x$ AND MINIMIZE W.R.T z :

$$\nabla \left[f(x) + \nabla f(x)^T z + \frac{M}{2} z^T I z \right] = 0 \Leftrightarrow$$

$$\nabla f(x) + Mz = 0 \Leftrightarrow z = -\frac{1}{M} \nabla f(x) \Leftrightarrow$$

$y = x - \frac{1}{M} \nabla f(x)$, SO MINIMUM R.H.S. IS

$$f(x) + \nabla f(x)^T \left(-\frac{1}{M} \nabla f(x) \right) + \frac{M}{2} \left(\frac{1}{M^2} \|\nabla f(x)\|_2^2 \right)$$

$$= f(x) - \frac{1}{2M} \|\nabla f(x)\|_2^2 \quad \text{IT FOLLOWS THAT } \forall y,$$

$$f(y) \leq f(x) - \frac{1}{2M} \|\nabla f(x)\|_2^2 \quad \text{AND THE RESULT FOLLOWS}$$

CONDITION NUMBER OF SUBLEVEL SETS:

$$\left. \begin{array}{l} (9.7) \\ (9.12) \end{array} \right\} \Rightarrow mI \leq \nabla^2 f(x) \leq MI \quad \forall x \in S$$

$$\Rightarrow K = \frac{M}{m} \text{ IS UPPER BOUND ON COND} \left(\nabla^2 f(x) \right) \triangleq \frac{\lambda_{\max}(\nabla^2 f(x))}{\lambda_{\min}(\nabla^2 f(x))}$$

K IS VERY IMPORTANT, BECAUSE IT AFFECTS THE SHAPE OF SUBLEVEL SETS, AND, THROUGH THEM, CONVERGENCE RATES OF ALGORITHMS. WE EXPLORE THE SHAPE MORE:

$$\text{LET } mI \leq \nabla^2 f(x) \leq MI \quad \forall x \in S$$

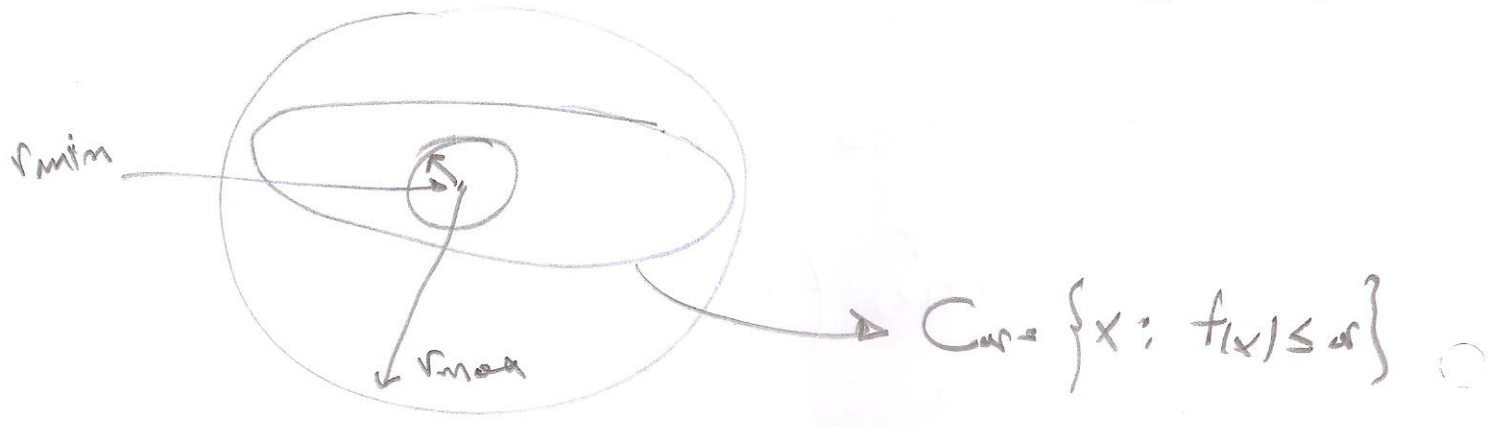
$$\left. \begin{array}{l} (9.13) \\ (9.8) \\ x = x^* \end{array} \right\} \Rightarrow p^* + \frac{M}{2} \|y-x\|_2^2 \geq f(y) \geq p^* + \left(\frac{m}{2} \right) \|y-x\|_2^2$$

SET $f(y) = \alpha$

$$p^* + \frac{M}{2} \|y - x^*\|_2^2 \geq \alpha \Rightarrow \|y - x^*\|_2 \geq r_{\min} = \left[\frac{2(\alpha - p)}{M} \right]^{\frac{1}{2}}$$

LOWER SET:

$$p^* + \left(\frac{m}{2}\right) \|y - x^*\|_2^2 \leq \alpha \Rightarrow \|y - x^*\|_2 \leq r_{\max} = \left[\frac{2(\alpha - p)}{m} \right]^{\frac{1}{2}}$$



SO WHEN $k \approx 1$ THE SUBLEVEL SET ARE GUARANTEED TO BE APPROXIMATELY CIRCULAR.

TAYLOR SERIES JUSTIFICATION

$$f(y) \approx p^* + \nabla f(x^*) (y - x^*) + \frac{1}{2} (y - x^*)^T \nabla^2 f(x^*) (y - x^*)$$

0 AT OPTIMUM x^*

WHEN $y \approx p^*$. THEREFORE, CLOSE TO OPTIMUM, EQUIPOTENTIAL LINES LOOK LIKE ELLIPSES WITH ECCENTRICITY DEPENDENT OF $\frac{M}{m}$

A LAST CAVEAT: M, m ARE UNKNOWN, MOST OF THE TIME. SO WHAT WE HAVE SEEN IS USEFUL FOR UNDERSTANDING PROBLEMS, AND NOT SO MUCH FOR SOLVING THEM.

9.2 DESCENT METHODS

ALL ALGORITHMS IN THIS CHAPTER LOOK LIKE THE FOLLOWING:

$$x^{(k+1)} = x^{(k)} + \underbrace{t^{(k)}}_{\text{STEP SIZE OR STEP LENGTH}} \underbrace{\Delta x^{(k)}}_{\in \mathbb{R}^n, \text{ SEARCH DIRECTION}}$$

ALL METHODS ARE DESCENT METHODS, I.E.

$$f(x^{(k+1)}) < f(x^{(k)}) \quad (\text{UNLESS } x^{(k)} \text{ IS OPTIMAL})$$

BY CONVEXITY WE HAVE
I.E. WE MOVE TOWARDS A DESCENT DIRECTION

$$\nabla f(x^{(k)})^T \Delta x^{(k)} < 0,$$

GENERAL DESCENT METHOD

GIVEN A STARTING POINT $x \in \text{dom} f$

REPEAT

- 1) FIND Δx
- 2) LINE SEARCH: CHOOSE $t > 0$
- 3) UPDATE $x := x + t \Delta x$

UNTIL STOPPING CRITERION (COMMON STOPPING CRITERION:

$$\|\nabla f(x)\|_2 \leq \epsilon)$$

EXACT LINE SEARCH:

$$t = \underset{s \geq 0}{\text{argmin}} \{f(x + s \Delta x)\}$$

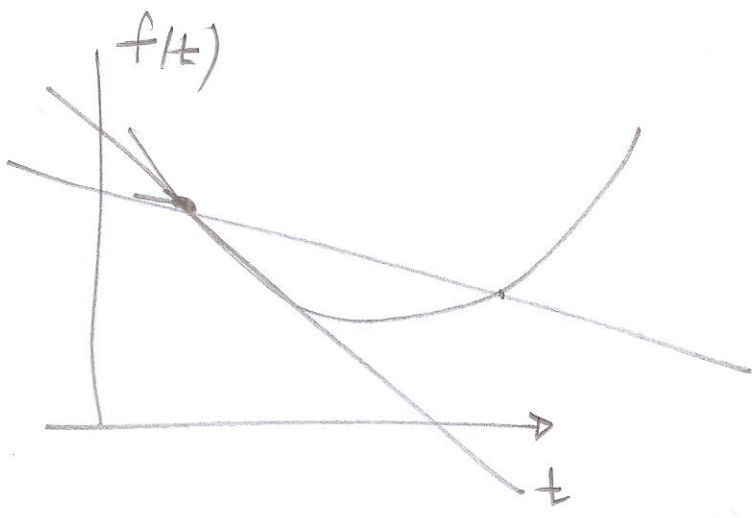
ONLY MAKES SENSE IF $f(x)$ IS EASY/FAST TO COMPUTE

BACKTRACKING LINE SEARCH:

GIVEN Δx FOR f AT $x \in \text{dom} f$, $\alpha \in (0, 0.5)$, $\beta \in (0, 1)$
 $t := 1$

WHILE $f(x + t \Delta x) > f(x) + \alpha t \nabla f(x)^T \Delta x$, $t := \beta t$

GEOMETRIC INTERPRETATION:



$$f(x + \Delta x) \approx f(x) + \nabla f(x)^T \Delta x < f(x) + \alpha \nabla f(x)^T \Delta x$$

THEREFORE, WE START AT DISTANCE Δx , AND MOVE CLOSER UNTIL WE GET A PERCENTAGE OF THE IMPROVEMENT WE DESIRE.

9.3 GRADIENT DESCENT METHOD

IT IS THE GENERAL DESCENT METHOD, WITH

$$\Delta x := -\nabla f(x)$$

CONVERGENCE ANALYSIS, EXACT LINE SEARCH

(I.E., HOW FAST WE FIND MINIMUM)

WE ASSUME: $mI \preceq \nabla^2 f(x) \preceq MI \quad \forall x \in S$

WE DEFINE $\tilde{f}(t) = f(x - t \nabla f(x))$

(9.13) $y = x - t \nabla f(x)$ } \Rightarrow

$$\tilde{f}(t) \leq f(x) - t \|\nabla f(x)\|_2^2 + \frac{Mt^2}{2} \|\nabla f(x)\|_2^2 \quad (9.17)$$

AT FIRST, LINEAR DECREASE, EVENTUALLY QUADRATIC DECREASE

MINIMIZE THE RIGHT HAND SIDE:

$$- \|\nabla f(x)\|_2^2 + M \epsilon \|\nabla f(x)\|_2^2 = 0 \Leftrightarrow \epsilon = \frac{1}{M}$$

AS WE STOP AT THE MINIMUM, WE HAVE

$$f(x^{(k+1)}) \leq f(x) - \frac{1}{M} \|\nabla f(x)\|_2^2 + \frac{1}{2M} \|\nabla f(x)\|_2^2$$

$$\Rightarrow f(x^{(k+1)}) - p^* \leq (f(x) - p^*) - \frac{1}{2M} \|\nabla f(x)\|_2^2$$

HOWEVER, WE KNOW THAT $\|\nabla f(x)\|_2^2 \geq 2m(f(x) - p^*)$
(CONVEX FROM (2.3))

$$\Rightarrow f(x^{(k+1)}) - p^* \leq \left(1 - \frac{m}{M}\right) (f(x) - p^*) \Rightarrow$$

$$f(x^{(k)}) - p^* \leq c^k (f(x^{(0)}) - p^*)$$
$$c = 1 - \frac{m}{M} < 1$$

SO RATE OF CONVERGENCE IS ^{BETTER} GEOMETRIC OR, AS WE SAY IN CONVERGENCE ANALYSIS, LINEAR

THEREFORE, IF WE WANT $f(x^{(k)}) - p^* \leq \epsilon$, WE NEED

$$c^k (f(x^{(0)}) - p^*) \leq \epsilon \Leftrightarrow$$

$$c^k \leq \frac{\epsilon}{f(x^{(0)}) - p^*} \Leftrightarrow k \log c \leq \log \left[\frac{\epsilon}{f(x^{(0)}) - p^*} \right]$$

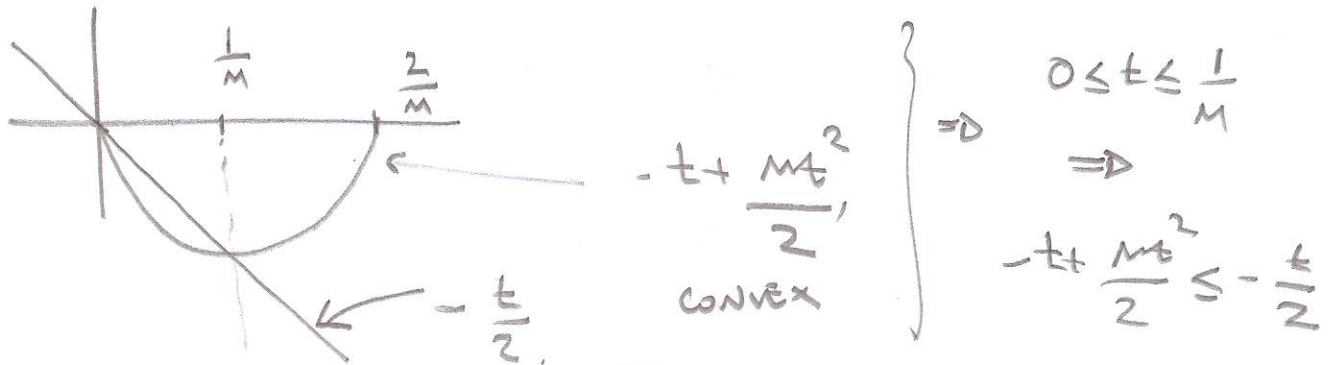
$$\Leftrightarrow k \geq \frac{\log \left[\frac{\epsilon}{f(x^{(0)}) - p^*} \right]}{\log c} \Leftrightarrow k \geq \frac{\log \left[\frac{f(x^{(0)}) - p^*}{\epsilon} \right]}{\log \left(\frac{1}{c} \right)}$$

THEREFORE, WE WANT $\frac{m}{M}$ AS LARGE AS POSSIBLE!

ANALYSIS FOR BACKTRACKING LINE SEARCH

(10)

WE SHOW THAT BACKTRACKING CONDITION IS SATISFIED FOR $0 < t \leq \frac{1}{M}$. THIS WILL GUARANTEE THAT WE DO NOT STOP TOO CLOSE TO THE ORIGIN $t=0$.



WE KNOW THAT

$$\tilde{f}(t) \leq f(x) - t \|\nabla f(x)\|_2 + \frac{Mt^2}{2} \|\nabla f(x)\|_2^2$$

(FOR $0 \leq t \leq \frac{1}{M}$)

$$\leq f(x) - \left(\frac{t}{2}\right) \|\nabla f(x)\|_2^2$$

$$\leq f(x) - \alpha t \|\nabla f(x)\|_2^2$$

(WE ASSUME $\alpha < \frac{1}{2}$)

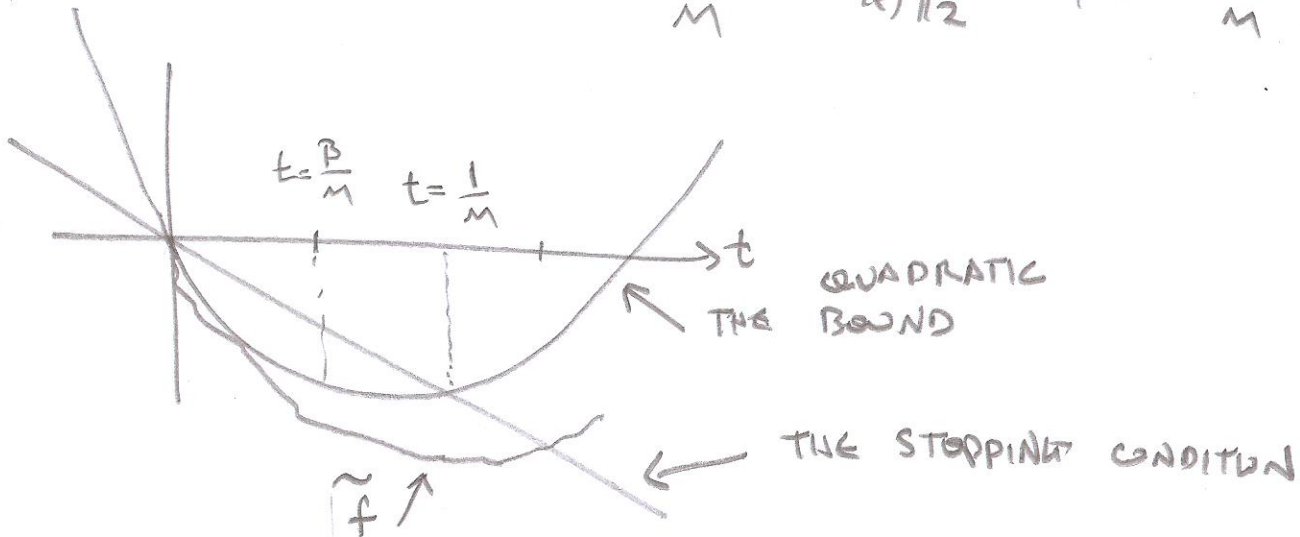
SO THE BACKTRACK IS GUARANTEED TO HAVE STOPPED EITHER

FOR SOME $t \geq \frac{\beta}{M}$, SO WE HAVE

FOR $t=1$, (IF $1 \leq \frac{1}{M}$)

$$f(x^+) \leq f(x) - \alpha \|\nabla f(x)\|_2^2 \quad (t=1)$$

$$\text{OR } f(x^+) \leq f(x) - \frac{\beta \alpha}{M} \|\nabla f(x)\|_2^2 \quad (t = \frac{\beta}{M})$$



(11)

$$\text{so } f(x^{+}) \leq f(x) - \min \left\{ \alpha, \frac{\beta \alpha}{M} \right\} \|\nabla f(x)\|_2^2$$

$$\Rightarrow f(x^{+}) - p^* \leq f(x) - p^* - \min \left\{ \alpha, \frac{\beta \alpha}{M} \right\} \|\nabla f(x)\|_2^2$$



$$\Rightarrow f(x^{(k+1)}) - p^* \leq (f(x^{(k)}) - p^*) \left[1 - 2m \min \left\{ \alpha, \frac{\beta \alpha}{M} \right\} \right] \geq 2m (f(x) - p^*)$$

"
C < 1

$$\Rightarrow f(x^{(k+1)}) - p^* \leq C^k [f(x^{(0)}) - p^*]$$

SO AGAIN WE HAVE LINEAR CONVERGENCE

AN EXAMPLE (QUADRATIC)

$$f(x) = \frac{1}{2} (x_1^2 + \gamma x_2^2)$$

$$H = \begin{bmatrix} 1 & 0 \\ 0 & \gamma \end{bmatrix}, \text{ so } m = \min(1, \gamma) \\ M = \max(1, \gamma).$$

IF WE START FROM $x^{(0)} = (\gamma, 1)$, WE CAN DERIVE CLOSED FORMULA EXPRESSIONS:

$$x_1^{(k)} = \gamma \left(\frac{\gamma-1}{\gamma+1} \right)^k, \quad x_2^{(k)} = \left(-\frac{\gamma-1}{\gamma+1} \right)^k$$

$$f(x^{(k)}) = \left(\frac{\gamma-1}{\gamma+1} \right)^{2k} f(x^{(0)})$$

THEREFORE, CONJUGATE IS INDEED LINEAR (EXCEPT $\gamma=1$)
WHEN $\gamma=1$, IN WHICH CASE WE FIND OPTIMUM IN FIRST TRY.

LET US COMPARE CONJUGATE WITH PREVIOUS BOUND,
(OF THIS EXAMPLE)

WE KNOW THAT, IN GENERAL CASE,
$$f(x^{(k)}) - p^* \leq C^k (f(x^{(0)}) - p^*)$$

WHERE $C = \left| 1 - \frac{m}{M} \right|$, $mI \leq \nabla^2 f(x) \leq MI$

IN THIS CASE,
$$f(x^{(k)}) - p^* \leq C^k (f(x^{(0)}) - p^*)$$

WHERE $C = \left(\frac{\gamma-1}{\gamma+1} \right)^2$

IF $\gamma > 1$, $\gamma = M$, $L = m \Rightarrow C = \left(\frac{M-m}{M+m} \right)^2 \Rightarrow$

$$C = \left(\frac{1 - \frac{m}{M}}{1 + \frac{m}{M}} \right)^2$$

IF $\gamma < 1$, $\gamma = m$, $M = L$, $\Rightarrow C = \left(\frac{m-M}{m+M} \right)^2 = \left(\frac{\frac{m}{M} - 1}{1 + \frac{m}{M}} \right)^2$

SO, ALWAYS,

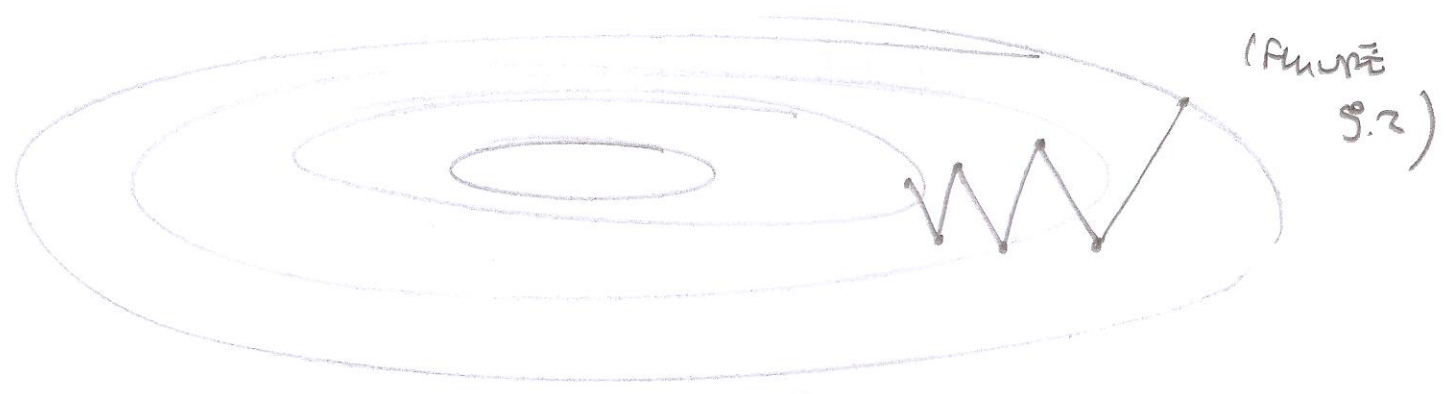
$$C = \left(\frac{1 - \frac{m}{M}}{1 + \frac{m}{M}} \right)^2$$

WHICH IS FASTER. IN FACT, FOR SMALL $\frac{m}{M}$, IT IS

4 TIMES FASTER, BECAUSE $C \approx \left(1 - \frac{m}{M} \right)^2 \left(1 - \frac{m}{M} \right)^2$

USING THE PROPERTY $\frac{1}{1+\epsilon} \approx 1-\epsilon$, $\epsilon \approx 0$.

SO INSTANT IN ANALYSIS IS ONLY BY A FACTOR OF 4.



THE ABOVE SHOWS WHY POOR CONDITION NUMBERS REALLY HURT US.

SOME GENERAL REMARKS:

- 1) CONVERGENCE IS LINEAR
- 2) α, β NOT SO CRUCIAL
- 3) WHEN CONDITION NUMBER IS HIGH, IE, $\mu \ll \kappa$, THE METHOD IS USELESS.

9.4 STEEPEST DESCENT METHODS

$$f(x+u) \approx \hat{f}(x+u) = f(x) + \underbrace{\nabla f(x)^T u}_{\substack{\uparrow \\ \text{DIRECTIONAL DERIVATIVE} \\ \text{IN DIRECTION } u}}$$

DEFINITION:

LET $\|\cdot\|$ BE ANY NORM. WE SELECT THE

NORMALIZED STEEPEST DESCENT DIRECTION

$$\begin{aligned} \Delta x_{msd} &= \arg \min \{ \nabla f(x)^T u \mid \|u\| = 1 \} \\ &= \arg \min \{ \nabla f(x)^T u \mid \|u\| \leq 1 \} \end{aligned}$$

DEFINITION: THE STEEPEST DESCENT METHOD IS LIKE THE GRADIENT DESCENT METHOD, BUT SELECTS AS THE DESCENT DIRECTION

$$\Delta x_{SD} = \|\nabla f(x)\|_* \Delta x_{USD}$$

UNNORMALIZED STEEPEST DESCENT

DUAL NORM, ENSURES THAT STEP IS ADEQUATELY AGGRESSIVE

VARIOUS CHOICES EXIST FOR THE CHOICE OF THE USED NORM:

① NORM IS EUCLIDEAN: THEN STEEPEST DESCENT = GRADIENT DESCENT

② NORM IS QUADRATIC NORM:

$$\|z\|_P = (z^T P z)^{1/2} = \|P^{1/2} z\|_2$$

WHERE $P \in S_{++}^m$. NOTE: AS $P \in S_{++}^m$, WE CAN ALWAYS WRITE

$$P = Q^T \Lambda Q$$

(ORTHOGONAL) $\text{diag}(\lambda_1, \lambda_2, \dots, \lambda_m)$, $\lambda_i > 0$,

SO
$$P^{1/2} \triangleq Q^T \text{diag}(\lambda_1^{1/2}, \lambda_2^{1/2}, \dots, \lambda_m^{1/2})$$

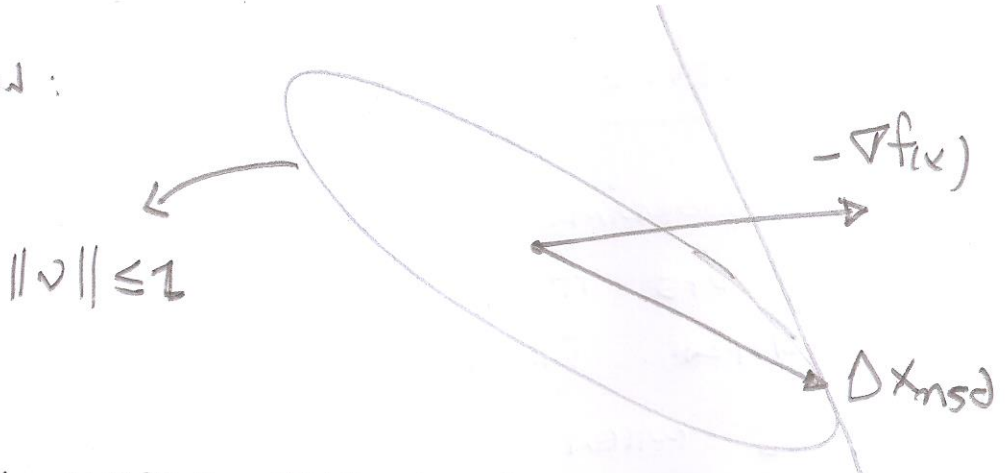
IN THIS CASE, WE CAN SHOW:

$$\Delta x_{msd} = -(\nabla f(x)^T P^{-1} \nabla f(x))^{-\frac{1}{2}} P^{-1} \nabla f(x),$$

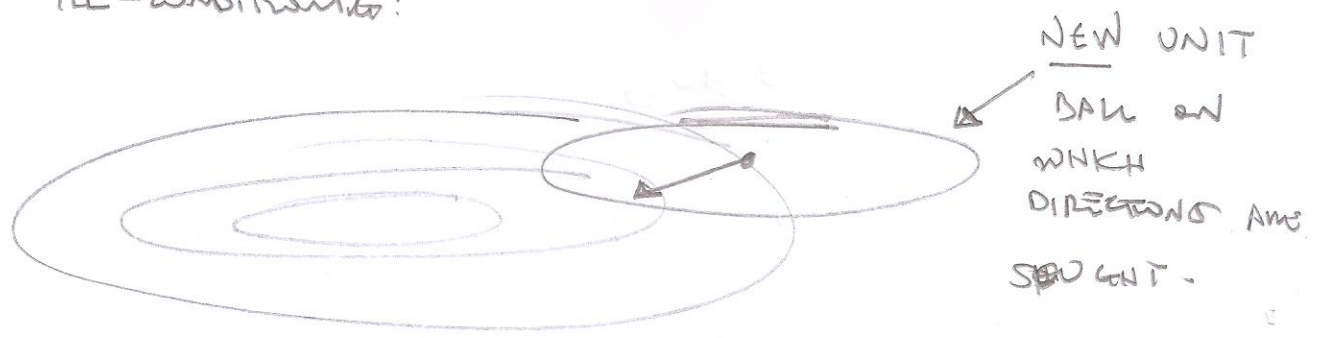
$$\Delta x_{sd} = -P^{-1} \nabla f(x),$$

BY SOLVING A CONSTRAINT OPTIMIZATION PROBLEM

INTUITIVELY:



SO, WHY WOULD THIS WORK BETTER? REMEMBER CASE OF ILL-CONDITIONING:



OFF COURSE, WE COULD DO MUCH MORE! IDEALLY, WE WANT TO MATCH WITH OUR UNIT BALL OUR HESSIAN AT THE OPTIMUM.

3) NORM IS l_1

$$\|z\|_1 = |z_1| + |z_2| + \dots + |z_n| \quad \text{IN THIS CASE,}$$

$$\Delta x_{msd} = \text{argmin} \left\{ \nabla f(x)^T v \mid \|v\|_1 \leq 1 \right\}.$$

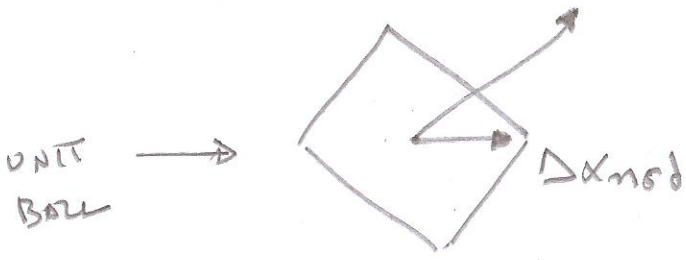
ONE OPTIMUM IS:
(OF THE POSSIBLY MANY)

$$\Delta x_{msd} = -\text{sign} \left(\frac{\partial f(x)}{\partial x_i} \right) e_i$$

norm of ∇ norm

$\Rightarrow \Delta x_{SD} = \Delta x_{NSD} \quad \|\nabla f(x)\|_{\infty} = -\frac{\partial f(x)}{\partial x_i} c_i$

INTUITION:



WE CAN ALWAYS PICK A DIRECTION IN A STANDARD BASIS VECTOR.

⇒ COORDINATE - DESCENT!

ALGORITHM: FIND COORDINATE ON WHICH DESCENT IS THE STEEPEST, AND USE IT TO DO EXACT LINE SEARCH OR BACKTRACK SEARCH. THEN REPEAT.

THIS IS VERY USEFUL WHEN SOLVING MINIMIZATION PROBLEM ACROSS ONE COORDINATE IS EASY.

9.5 NEWTON'S METHOD

DEFINITION: NEWTON STEP: $\Delta x_{nt} = -[\nabla^2 f(x)]^{-1} \nabla f(x)$.

IT IS CERTAINLY A DESCENT DIRECTION, BECAUSE

$\nabla f(x)^T \Delta x_{nt} = -\nabla f(x)^T [\nabla^2 f(x)]^{-1} \nabla f(x) < 0$,

SINCE $[\nabla^2 f(x)]^{-1}$ IS POSITIVE DEFINITE. THIS IS THE STEP WE WILL USE FOR OUR GRADIENT DESCENT ALGORITHM. WHY?

1) IT IS THE MINIMIZER OF THE SECOND ORDER APPROXIMATION

$\hat{f}(x+v) = f(x) + \nabla f(x)^T v + \frac{1}{2} v^T \nabla^2 f(x) v$.

SO LET'S TRY TO MINIMIZE THAT!

$\nabla \hat{f}(x+v) = \nabla f(x) + \nabla^2 f(x) v \Rightarrow v = -[\nabla^2 f(x)]^{-1} \nabla f(x)$

2) IT IS THE STEEPEST DESCENT DIRECTION UNDER THE HESSIAN NORM (WHICH WE KNOW IS A VERY GOOD CHOICE FOR THE S.D. NORM). INDEED:

FROM STEEPEST DESCENT THEORY, IF WE USE $\|z\|_p = (z^T P z)^{\frac{1}{2}} = \|P^{\frac{1}{2}} z\|_2$, THEN $\Delta x_{SD} = -P^{-1} \nabla f(x)$

3) IT IS THE SOLUTION OF THE LINEARIZED OPTIMALITY CONDITIONS
OBSERVE THAT

$$\nabla f(x+w) \approx \nabla f(x) + \nabla^2 f(x) w$$

INDEED, LET $F(x) = \begin{bmatrix} f_1(x) \\ \vdots \\ f_m(x) \end{bmatrix}$. THEN

$$F(x+w) \approx F(x) + \begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \frac{\partial f_1}{\partial x_2} & \dots & \frac{\partial f_1}{\partial x_n} \\ \frac{\partial f_2}{\partial x_1} & \frac{\partial f_2}{\partial x_2} & \dots & \frac{\partial f_2}{\partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \frac{\partial f_m}{\partial x_2} & \dots & \frac{\partial f_m}{\partial x_n} \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_n \end{bmatrix}$$

AND SUBSTITUTING $F(x)$ FOR ∇f , THE ABOVE RESULT FOLLOWS

SO MINIMIZING THE ABOVE, THE RESULT FOLLOWS.
IN OTHER WORDS, THIS IS AN ESTIMATION OF WHERE THE OPTIMALITY CONDITIONS HOLD

A BASIC PROPERTY OF THE NEWTON INCREMENT: IT IS INVARIANT TO AFFINE TRANSFORMATIONS.

LET $T \in \mathbb{R}^{m \times m}$, $\det T \neq 0$.

LET $\bar{f}(y) = f(Ty) \Rightarrow \begin{cases} \nabla \bar{f}(y) = T^T \nabla f(x) \\ \nabla^2 \bar{f}(y) = T^T \nabla^2 f(x) T \end{cases}$

$x = Ty$

NEWTON STEP FOR \bar{f} AT y IS:

$$\begin{aligned} \Delta y_{nt} &= - \left(T^T \nabla^2 f(x) T \right)^{-1} \left(T^T \nabla f(x) \right) \\ &= - T^{-1} \left[\nabla^2 f(x) \right]^{-1} \nabla f(x) = T^{-1} \Delta x_{nt} \Rightarrow \end{aligned}$$

$$\Delta x_{nt} = T \Delta y_{nt}$$

\Rightarrow

$$x + \Delta x_{nt} = T (y + \Delta y_{nt})$$

THIS IMPLIES THAT NEWTON'S METHOD, TO BE INTRODUCED BELOW, IS NOT AFFECTED BY COORDINATE TRANSFORMATIONS.

ALGORITHM: NEWTON'S METHOD

GIVEN α STARTING POINT AND TOLERANCE $\epsilon > 0$.

REPEAT

1) COMPUTE $\left\{ \begin{array}{l} \text{NEWTON STEP } \Delta x_{nt} := - \nabla^2 f(x)^{-1} \nabla f(x) \\ \text{NEWTON DECREMENT } \lambda := \nabla f(x)^T \nabla^2 f(x) \nabla f(x) \end{array} \right.$

2) STOPPING CRITERION: EXIT IF $\frac{\lambda^2}{2} \leq \epsilon$

3) LINE SEARCH: CHOOSE STEP SIZE t BY BACKTRACKING

4) UPDATE $x := x + t \Delta x_{nt}$

COMMENTS:

1) IF THERE IS NO BACKTRACK, AND $t=1$ ALWAYS, THEN WE CALL THE METHOD THE PURE NEWTON'S METHOD.

2) REGARDING THE STOPPING CRITERION:

$$f(x) - \inf_y \hat{f}(y) = f(x) - \hat{f}(x + \Delta x_{nt}) =$$

$$-\nabla f(x)^T \left[-\nabla^2 f(x)^{-1} \right] \nabla f(x)$$

THE QUADRATIC APPROXIMATION OF f

$$- \frac{1}{2} \nabla f(x)^T \left[-\nabla^2 f(x)^{-1} \right]^T \nabla^2 f(x) \left[-\nabla^2 f(x)^{-1} \nabla f(x) \right]$$

$$= + \frac{1}{2} \nabla f(x)^T \left[\nabla^2 f(x)^{-1} \right] \nabla f(x) = \frac{\lambda^2}{2}$$

SO $\frac{\lambda^2}{2}$ IS AN ESTIMATE OF OUR ERROR.

3) ONCE WE ARE CLOSE ENOUGH TO THE SOLUTION, THE CONVERGENCE IS QUADRATIC!

ADVANTAGES OF NEWTON'S METHOD

- 1) RAPID CONVERGENCE
- 2) AFFINE INVARIANT
- 3) SCALES WELL
- 4) DOES NOT DEPEND ON CHOICE OF COORDINATES

DISADVANTAGE: YOU NEED TO FIND A MATRIX (THE HESSIAN) AND INVERT IT!

⇒ QUASI-NEWTON METHODS HAVE BEEN COMPUTATIONAL DEVELOPED, WHICH AVOID HEAVY